

Project 3: Analyses of daily COVID-19 cases across nations

Jiayi Shen (js5354), Siquan Wang (sw3442), Jack Yan (xy2395)

5/1/2020

1. Introduction

The pandemic of COVID-19 is the biggest challenge that the world is facing right now. Our lives are all deeply affected by this public health crisis. By building a model on the growth of COVID-19 cases, we can have a better understanding of the current status and then plan future responses. Thus analyzing existing data and predicting future trajectories has become the most important task faced by public health expertises and policy makers.

The objective of this project includings:

- Develop an optimization algorithm to fit a logisitic curve to each region, using the global COVID-19 data;
- Evaluate how appropriate it is to use such model;
- Apply clustering methods to the logistic grwoth model parameters, and observe the patterns.

2. Data

The datasets used in this project are adapted from https://github.com/CSSEGISandData/COVID-19/tree/master/csse_covid_19_data/csse_covid_19_daily_reports. We used the latest data updated on 04/29/2020, so that we have more data points. The data recorded the following variables:

Id: Record ID

Province/State: The lcoal state/province of the record; 54% records do not have this info;

Country/Region: The country/regionoof the record;

Lat: Lattudiute of the record;

Long: Longitude of the record;

Date: Date of the record; from Jan 21 to March 23;

ConfirmedCases: The number of confirme case on that day;

Fatalities: The number of death on that day;

For the purpose of fitting logistic growth curve, the main variables of interest are **Country/Region** and **ConfirmedCases**. We groupped the dataset by **Country/Region**, and pulled out **days since first case** and **cumulative confirmed cases on each day**.

3. Method

3.1 Logisitic curves

Logisitic curves could be one way to model the trajectory of cumulative cases; It is a parametric function with the form

$$f(t) = \frac{a}{1 + \exp\{-b(t - c)\}},$$

where t is the days since the first infection; a is the upper bound, i.e. the maximum number of cases a region can reach, b is growth rate, and c is the mid-point, where the curve changes from convex to concave; Each curve is uniquely defined by (a, b, c) . By design a logistic curve increases exponentially at beginning and slows down at the end.

Consider the most commonly used loss function: Mean Squared Error (MSE), which takes the form

$$l = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n \left(Y_i - \frac{a}{1 + \exp\{-b(T_i - c)\}} \right)^2$$

Then our goal is to minimize the above loss function with respect to (a, b, c) . To do this, we can apply Newton-Raphson algorithm. Let $\theta = (a, b, c)$. Newton's method suggests to update θ iteratively, such that the i th step is given by

$$\theta_i = \theta_{i-1} - [\nabla^2 l(\theta_{i-1})]^{-1} \nabla l(\theta_{i-1})$$

where $\nabla l(\theta_{i-1})$ is the gradient, and $[\nabla^2 l(\theta_{i-1})]^{-1}$ is the Hessian matrix.

In this particular case, we replace the Hessian matrix with an identity matrix to simplify the computation and increase the efficiency. Step-halving is also incorporated to control the step size. Then the i th step of our Newton algorithm is given by

$$\theta_i = \theta_{i-1} + \lambda \mathbf{H}_{i-1, p \times p} \nabla l(\theta_{i-1})$$

where $\mathbf{H}_{i-1, p \times p} = I_{p \times p}$, $\lambda \in (0, 1)$.

The gradient vector $\nabla l(\theta_{i-1})$ of our loss function is:

$$\begin{pmatrix} \partial l / \partial a \\ \partial l / \partial b \\ \partial l / \partial c \end{pmatrix} = \begin{pmatrix} \frac{2}{n} \sum_{i=1}^n \left(Y_i - \frac{a}{1 + \exp\{-b(T_i - c)\}} \right) \cdot \frac{-1}{1 + \exp\{-b(T_i - c)\}} \\ \frac{2}{n} \sum_{i=1}^n \left(Y_i - \frac{a}{1 + \exp\{-b(T_i - c)\}} \right) \cdot \frac{-a(T_i - c) \exp\{-b(T_i - c)\}}{(1 + \exp\{-b(T_i - c)\})^2} \\ \frac{2}{n} \sum_{i=1}^n \left(Y_i - \frac{a}{1 + \exp\{-b(T_i - c)\}} \right) \cdot \frac{ab \exp\{-b(T_i - c)\}}{(1 + \exp\{-b(T_i - c)\})^2} \end{pmatrix} \quad (1)$$

We set the convergence criteria to be that the difference in MSE of two consecutive iterations is smaller than 10^{-5} , and the maximum iteration number to be 1000.

3.2 Clustering

To understand which countries/regions are similar in terms of the trajectory of COVID-19 cases, we applied two clustering methods, K-mean and Gaussian mixture model, to group the fitted parameters $(\hat{a}, \hat{b}, \hat{c})$.

The gaussian mixture model assumes that $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \in \mathbb{R}^p$ are i.i.d. random vectors following a mixture multivariate normal distributions with k hidden groups. In this case, $\mathbf{x}_i = (Y_i, T_i)'$ and $(a, b, c) \in \mathbb{R}^3$. And

$$\mathbf{x}_i \sim \begin{cases} N(\boldsymbol{\mu}_1, \Sigma_1), \text{ with probability } p_1 \\ N(\boldsymbol{\mu}_2, \Sigma_2), \text{ with probability } p_2 \\ \vdots, \quad \quad \quad \vdots \\ N(\boldsymbol{\mu}_k, \Sigma_k), \text{ with probability } p_k \end{cases}$$

$$\mathbf{x}_i \sim \begin{cases} N(\boldsymbol{\mu}_1, \Sigma_1), \text{ with probability } p_1 \\ N(\boldsymbol{\mu}_2, \Sigma_2), \text{ with probability } p_2 \\ \vdots, \quad \quad \quad \vdots \\ N(\boldsymbol{\mu}_k, \Sigma_k), \text{ with probability } p_k \end{cases}$$

The density of a multivariate normal \mathbf{x}_i is

$$f(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \frac{\exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}))}{\sqrt{(2\pi)^P |\Sigma|}}$$

The observed likelihood of $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ is

$$L(\theta; \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = \prod_{i=1}^n \sum_{j=1}^k p_j f(\mathbf{x}_i; \boldsymbol{\mu}_j, \Sigma_j)$$

Let $\mathbf{r}_i = (r_{i,1}, \dots, r_{i,k}) \in \mathbb{R}^k$ as the cluster indicator of \mathbf{x}_i , which takes form $(0, 0, \dots, 0, 1, 0, 0)$ with $r_{i,j} = I\{\mathbf{x}_i \text{ belongs to cluster } j\}$. The cluster indicator \mathbf{r}_i is a latent variable that cannot be observed. Therefore, we use EM algorithm to iteratively estimate model parameters.

E-step: Evaluate the responsibilities using the current parameter values

$$\gamma_{i,k}^{(t)} = P(r_{i,k} = 1 | \mathbf{x}_i, \theta^{(t)}) = \frac{p_k^{(t)} f(\mathbf{x}_i | \boldsymbol{\mu}_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K p_j^{(t)} f(\mathbf{x}_i | \boldsymbol{\mu}_j^{(t)}, \Sigma_j^{(t)})}$$

M-step:

$$\theta^{(t+1)} = \arg \max \ell(\mathbf{x}, \gamma^{(t)}, \theta).$$

Let $n_k = \sum_{i=1}^n \gamma_{i,k}$, we have

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{1}{n_k} \sum_{i=1}^n \gamma_{i,k} \mathbf{x}_i$$

$$\Sigma_k^{(t+1)} = \frac{1}{n_k} \sum_{i=1}^n \gamma_{i,k} (\mathbf{x}_i - \boldsymbol{\mu}_k^{(t+1)})(\mathbf{x}_i - \boldsymbol{\mu}_k^{(t+1)})^T$$

$$p_k^{(t+1)} = \frac{n_k}{n}$$

By iteratively updating the E-step and the M-step, we can reach a solution upon convergence.

On the other hand, the K -means algorithm essentially finds cluster centers and cluster assignments that minimize the objective function

$$J(\mathbf{r}, \boldsymbol{\mu}) = \sum_{i=1}^n \sum_{j=1}^k r_{i,j} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2$$

where $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k\}$ are the centers of the k (unknown) clusters, and $\mathbf{r}_i = (r_{i,1}, \dots, r_{i,k}) \in \mathbb{R}^k$ as the *hard* cluster assignment of \mathbf{x}_i .

4. Results

4.1 Logistic curve

Logistic curves were used to fit cumulative confirmed cases in each country/region. The fitted results for some selected countries are shown in **Figure 1**. According to the model, as of 04/29/2020, 37 regions have passed the midpoint, among which 100 regions are close to the end of virus spreading. If the confirmed cases in a region is greater than 95% of the estimated \mathbf{a} (upper bound) in the logistic curve, we define the region as close to the end of virus spreading.

In general, Regions with small population size or at an early stage of spreading tend to have a more flat rate of growth, such as Bhutan and Venezuela (**Figure 2**).

Task 1.2

Do you think if the logistic curve is a reasonable model for fitting the curmulative casesa and predicting future new cases?

I think logistic curve might not be an appropriate model for fitting the curmulative cases and predicting future new cases. First our Newton's algorithm suggests that although we spent much efforts in optimizing this optimization algorithm (including standardization, replacing Hessian by Identity matrix and step-halving), the convergence speed might still have some unknown internal issue and the results highly depending on the initial starting value if we do not do standardization, which means that our model is not that robust and the estimated curve shape parameters might be questionable in some unusual case. Second, some African countries' results are little bit strange base on our model, which might be due to too little data available. Finally, I may suggest using some other exponential-distribution based curves to fit the data and incorporate more shape paramters in the model fitting process. Or we could some hybrid approaches to fit different types of curves based on preclassification of countries.

Task 2: Clustering your fitted curves

Apply K-mean and Guassian mixture model (with EM algorithm) to cluster the fitted parameters $(\hat{a}, \hat{b}, \hat{c})$; Which algoirhtm does a better job in clustering those curves? Are the resulting clusters related to geogrpahic regions, or the starting timeing of the local virus spreading, or the resources of the regions? You may use external informations to help understand the clusters, i.e. find plausible explanations why some regions have similar (a, b, c) ?

Figures

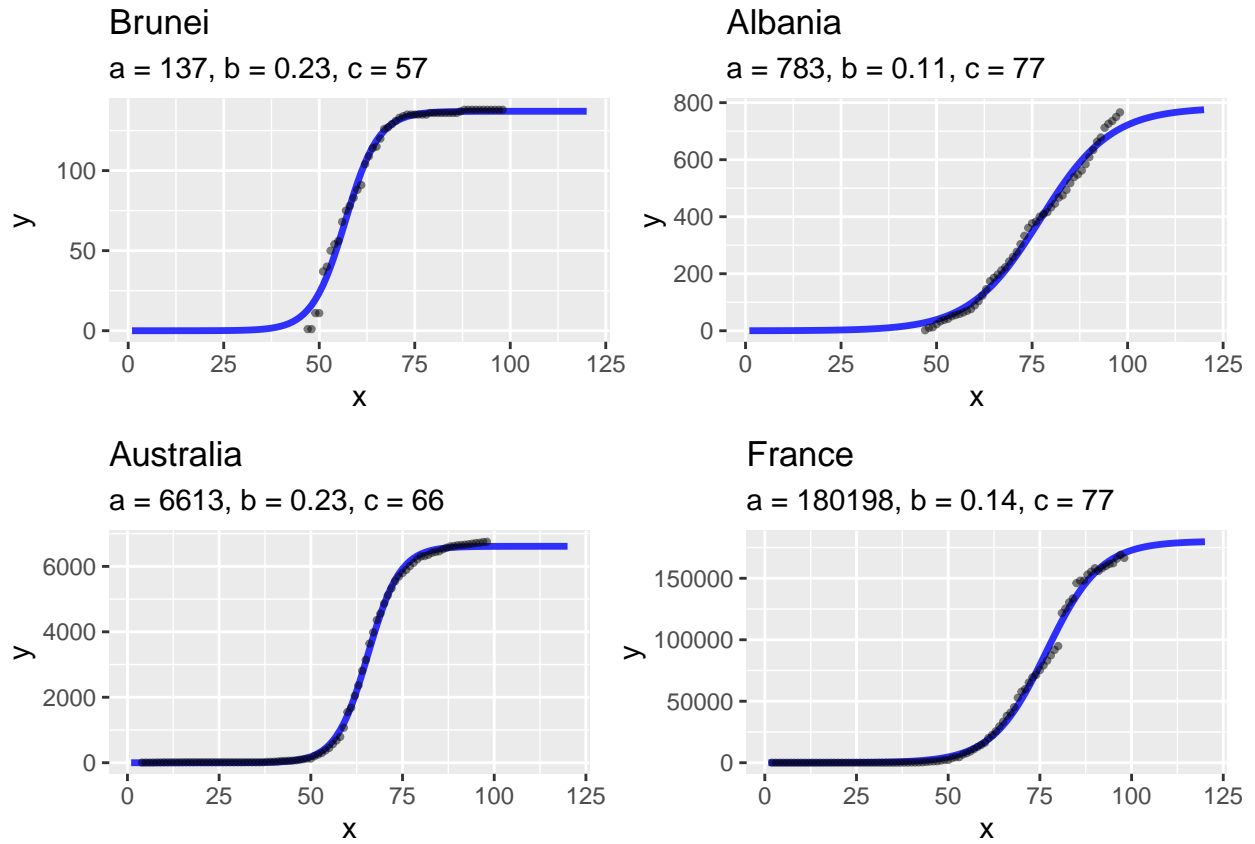


Figure 1 Logistic curves fitted on cumulative confirmed cases in selected countries

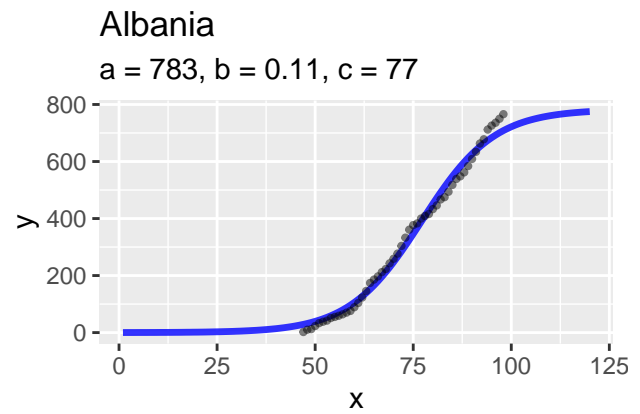
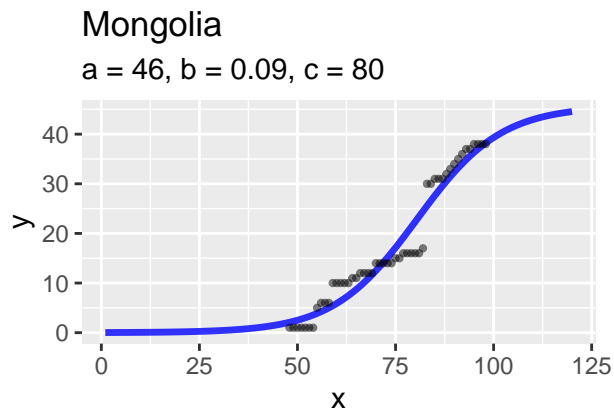
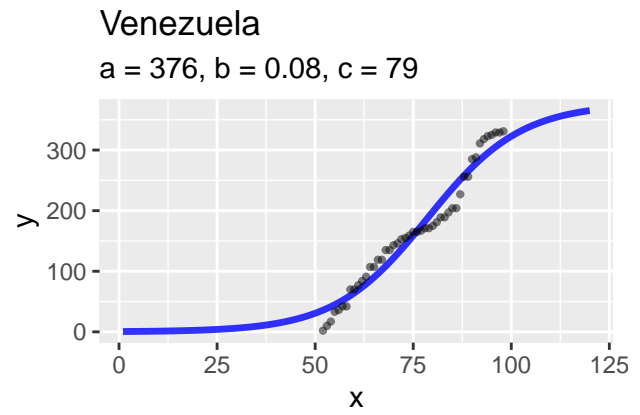
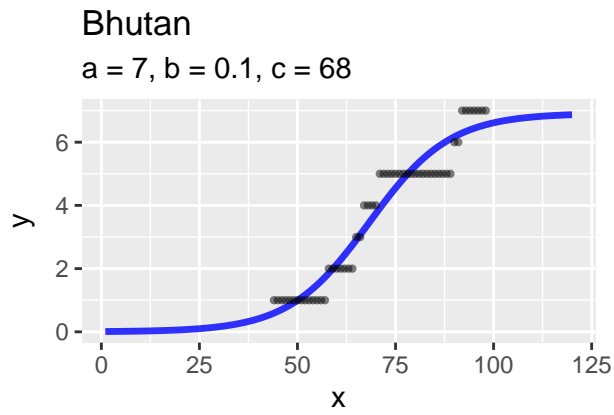


Figure 2 Logistic curves of selected countries with flat rate of growth

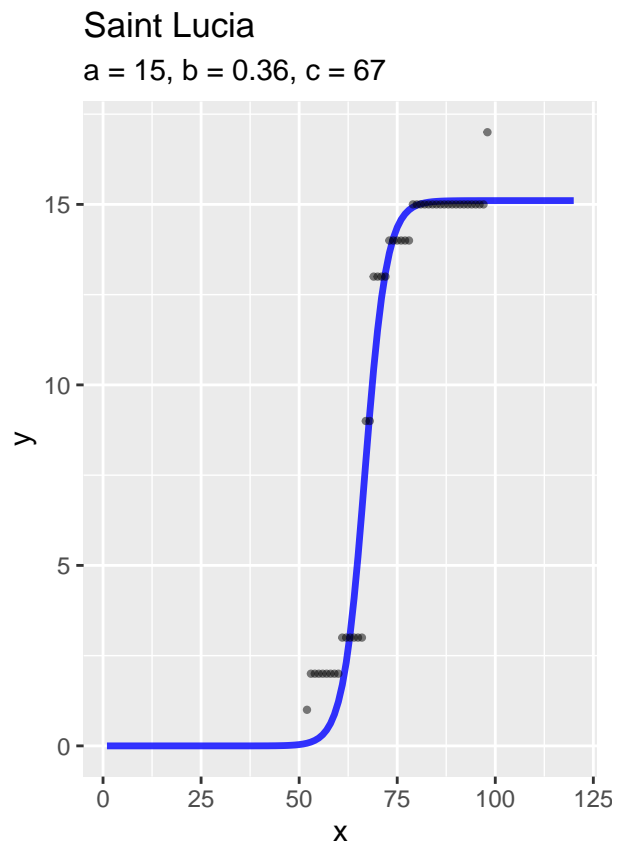
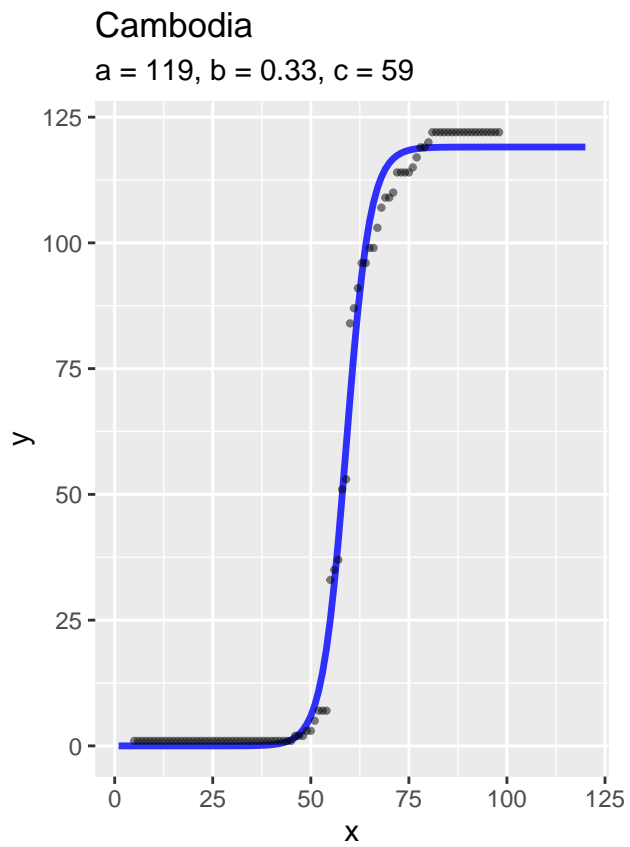


Figure 3 Logistic curves of selected countries with steep rate of growth